

In this problem sheet, we consider the problem of linear regression with p predictors and one intercept,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

where $\mathbf{y}^t = (y_1, \dots, y_n)$ is the column vector of target values, $\boldsymbol{\beta}^t = (\beta_0, \dots, \beta_p)$ is the column vector of coefficients, $\boldsymbol{\epsilon}^t = (\epsilon_1, \dots, \epsilon_n)$ is the vector of random errors, and \mathbf{X} is the $n \times (p + 1)$ matrix of observations given by

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix}.$$

We assume that $\boldsymbol{\epsilon}$ has a multivariate normal distribution with covariance matrix $\sigma^2 I$. The case $p = 1$ is referred to as simple linear regression.

Problem 0.

- (i) Show that the least square solution is given by $\hat{\boldsymbol{\beta}} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{y}$. Argue that if \mathbf{X} has rank $p + 1$, then $\mathbf{X}^t \mathbf{X}$ is indeed invertible.
- (ii) Let $\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{H} \mathbf{y}$, where \mathbf{H} is referred to as the hat matrix. Denote by \mathbf{x}_j the j -th column of \mathbf{X} . Show that the inner product between $\hat{\mathbf{y}} - \mathbf{y}$ and \mathbf{x}_j is 0. Deduce a geometrical interpretation of the least square solution.
- (iii) Explain with words the interpretation of the regression coefficient $\hat{\beta}_j$ in a multivariate linear regression setting, with non-orthogonal inputs.
- (iv) Let $\mathbf{X} = \mathbf{Q} \mathbf{R}$ be the QR decomposition of \mathbf{X} . Give a geometrical interpretation of the matrix \mathbf{Q} . Then show that $\hat{\mathbf{y}} = \mathbf{Q} \mathbf{Q}^t \mathbf{y}$.
- (v) Show that $\hat{\boldsymbol{\beta}}$ is an unbiased estimator of $\boldsymbol{\beta}$, and derive the covariance matrix of $\hat{\boldsymbol{\beta}}$.
- (vi) Show that $\hat{\sigma}^2 = \sum_i (y_i - \hat{y}_i)^2 / (n - p - 1)$ is unbiased for σ^2 .
- (vii) How would you test if a particular predictor is associated with the response variable? Under normally distributed errors $\boldsymbol{\epsilon}$, which test would you use?
- (viii) What is the difference between a confidence interval and a prediction interval? Which one is wider and why?
- (ix) Give a hypothesis test to detect outliers.
- (x) How would you check other model assumptions such as normality and constant variance?

Problem 1.

The least square estimator $\hat{\beta}$ of β is optimal in a certain sense. This is made precise with the result below, known as the Gauss-Markov Theorem.

Gauss-Markov Theorem. The least square estimator $\hat{\beta}$ has minimum variance amongst all unbiased linear estimators of β .

Two remarks:

- (i) Linear should be understood as linear with respect to \mathbf{y} , that is of the form $\mathbf{B}\mathbf{y}$, where \mathbf{B} is some $(p+1) \times n$ matrix. This way, the LS estimator $\hat{\beta}$ is indeed linear since $\hat{\beta} = (\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t\mathbf{y}$.
- (ii) It is possible to define a partial ordering in the class of real symmetric matrices. We say that $\mathbf{B}_1 \preceq \mathbf{B}_2$ if $\mathbf{B}_2 - \mathbf{B}_1$ is a positive semi-definite matrix. In other words, we have that for any vector \mathbf{x} , $\mathbf{x}^t\mathbf{B}_1\mathbf{x} \leq \mathbf{x}^t\mathbf{B}_2\mathbf{x}$. Equivalently, the matrix $\mathbf{B}_2 - \mathbf{B}_1$ has non-negative eigenvalues.

The goal of this problem is to prove the Gauss-Markov Theorem.

- (a) Let $\tilde{\beta} = \mathbf{B}\mathbf{y}$ be another unbiased linear estimator of β . Show that $\mathbf{B}\mathbf{X} = \mathbf{I}$.
- (b) Show that $\text{Cov}(\tilde{\beta} - \hat{\beta}, \hat{\beta}) = 0$.
- (c) Show that for two random vectors U and V , the covariance matrix Σ_{U+V} of the random vector $U + V$ satisfies

$$\Sigma_{U+V} = \Sigma_U + \Sigma_V + \text{Cov}(U, V) + \text{Cov}(V, U).$$

- (d) Deduce from (b) and (c) that

$$\Sigma_{\tilde{\beta}} = \Sigma_{\tilde{\beta} - \hat{\beta}} + \Sigma_{\hat{\beta}},$$

and conclude.

Problem 2.

Using geometrical considerations and Pythagora's theorem, we saw during the lectures that for the general linear regression model with intercept, the Total Sum of Squares (TSS) can be decomposed as a sum of Explained Sum of Squares (ESS) and Residual Sum of Squares (RSS),

$$\begin{aligned} TSS &= \sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{\hat{y}})^2 + \sum (y_i - \hat{y}_i)^2 \\ &= ESS + RSS, \end{aligned}$$

where $\bar{y} = n^{-1} \sum y_i$ and $\bar{\hat{y}} = n^{-1} \sum \hat{y}_i$.

- (a) Check that this decomposition holds using direct calculations.

Hint: You may use the fact that $\mathbf{X}^t(\mathbf{y} - \mathbf{X}\hat{\beta}) = 0$, and conclude from the first equation corresponding to the column of ones in \mathbf{X} that $\bar{y} = \bar{\hat{y}}$.

(b) Show that in the case of simple linear regression, the R^2 coefficient defined as

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}$$

is equal to the square of the empirical correlation coefficient r , defined as

$$r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{(\sum(x_i - \bar{x})^2)(\sum(y_i - \bar{y})^2)}}.$$

Problem 3.

Show that for the simple linear regression model, the variance of the LS coefficient estimates are given by

$$\text{Var } \hat{\beta}_0 = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}}{\sum(x_i - \bar{x})^2} \right) \quad \text{Var } \hat{\beta}_1 = \frac{\sigma^2}{\sum(x_i - \bar{x})^2}.$$

Note that from the expression of $\text{Var } \hat{\beta}_1$, the more variability there is in the x s, and the less there is in the estimation of the slope: we get a better estimate as the input variable is more spread out. This makes sense, right?

Problem 4.

In linear regression, the k -th diagonal element h_{kk} of the projection matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}^t\mathbf{X})^{-1}\mathbf{X}^t$, referred to as the leverage of x_k , quantifies how much observation k contributes to the LS estimate.

(a) Show that for simple linear regression, the leverage corresponding to the k -th observation can be written

$$h_{kk} = \frac{1}{n} + \frac{(x_k - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2},$$

Conclude that necessarily, h_{kk} belongs to the interval $[1/n, 1]$. Which observation corresponds to the smallest leverage?

(b) Show that for a multivariate regression model, $0 \leq h_{kk} \leq 1$.

Hint: Use that fact that \mathbf{H} is idempotent and symmetrical.

(c) If $h_{kk} = 0$ or 1 , then $h_{kj} = 0$ for all $j \neq k$.

(d) For all $j \neq k$, $-1/2 \leq h_{kj} \leq 1/2$.

Problem 5.

Consider multivariate linear regression with independent errors $\epsilon \sim \mathcal{N}(0, \sigma^2)$.

- (a) Show that the least square estimator of β is equal to the maximum likelihood estimator.
- (b) What is the maximum likelihood estimator of σ^2 ? Compare its expression with the unbiased estimator

$$\hat{\sigma}^2 = \frac{1}{n-p-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$

Problem 6.

Prove the Sherman-Morrison-Woodbury Theorem: for any non-singular $p \times p$ matrix \mathbf{A} and $p \times 1$ column vectors \mathbf{u} and \mathbf{v} ,

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^t)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1}\mathbf{u}\mathbf{v}^t\mathbf{A}^{-1}}{1 + \mathbf{v}^t\mathbf{A}^{-1}\mathbf{u}}.$$

Problem 7.

We consider simple linear regression, with one predictor and an intercept. The LS estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ are computed from a training sample $\{(x_i, y_i)\}$ of size n .

- (i) Show that $\hat{\beta}_0 + \hat{\beta}_1 x \sim \mathcal{N}(\beta_0 + \beta_1 x, \gamma_n \sigma^2)$, where

$$\gamma_n = \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum (x_i - \bar{x})^2}.$$

- (ii) Deduce from (i) a confidence interval for the mean response at x .
- (iii) How would you modify the confidence interval in (ii) to obtain a prediction interval for the response variable y at x ?

Problem 8.

Consider two linear regression models

$$\text{Model 1: } Y = \beta_0 + \beta_1 X + \epsilon,$$

$$\text{Model 0: } Y = \beta_0 + \epsilon,$$

We want to test if the slope is needed in our model (simpler Model 0 preferred over Model 1). That is, we want to test for the null hypothesis $H_0 : \beta_1 = 0$. For this we have two options: a t-test or an F-test. The goal of this problem is to show that in this simple setting, the two tests are equal.

- (i) t-test. Let \mathbf{X} be the $(n \times 2)$ the matrix of observations, the first column being a column of ones, and $\mathbf{y}^t = (y_1, \dots, y_n)$ be the column vector of outputs in the training data set. Consider Model 1,

$$\mathbf{y} = \mathbf{X}\beta + \epsilon,$$

where $\beta^t = (\beta_0, \beta_1)$, and ϵ has a multivariate normal distribution with covariance matrix $\sigma^2 I$. We denote by $\hat{\beta}^t = (\hat{\beta}_0, \hat{\beta}_1)$ the least square estimate of β . Let v_{ij} be the entry in the i -th row and j -th column of the 2×2 matrix $(\mathbf{X}^t \mathbf{X})^{-1}$. Show that

$$z := \frac{\hat{\beta}_1}{\hat{\sigma} \sqrt{v_{22}}} \sim t_{n-2},$$

where $\hat{\sigma}^2$ is an unbiased estimator of σ^2

(ii) F-test. Denote by RSS_1 (respectively RSS_0) the residual sum of squares of the larger model (respectively, of the reduced model). We saw during the lectures that

$$F = \frac{(RSS_0 - RSS_1)}{RSS_1 / (n - 2)} \sim F_{1, n-2}.$$

Show that we can re-express F as

$$F = \frac{\hat{\beta}_1^2}{\hat{\sigma}^2} \sum (x_i - \bar{x})^2.$$

(iii) Conclude that the t-test and F-test are equal.

Problem 9.

One of the goals of linear regression is to predict the value of a target variable y for a new observation x . Let $\mathbf{x}_{n+1}^t = (1, x_{n+1,1}, \dots, x_{n+1,p})$ be a new observation. We model the response variable by

$$y_{n+1} = \mathbf{x}_{n+1}^t \beta + \epsilon_{n+1},$$

where $\mathbf{E} \epsilon_{n+1} = 0$, $\text{Var} \epsilon_{n+1} = \sigma^2$, and $\text{Cov}(\epsilon_{n+1}, \epsilon_i) = 0$, for $i = 1, \dots, n$. The target value y_{n+1} is predicted using $\hat{y}_{n+1} = \mathbf{x}_{n+1}^t \hat{\beta}$, where $\hat{\beta} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{y}$ is the least square estimate of β computed from the training data. The prediction error is defined as $\hat{\epsilon}_{n+1} = y_{n+1} - \hat{y}_{n+1}$.

(a) Show that

$$\begin{aligned} \mathbf{E} \hat{\epsilon}_{n+1} &= 0, \\ \text{Var} \hat{\epsilon}_{n+1} &= \sigma^2 \left(1 + \mathbf{x}_{n+1}^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}_{n+1} \right). \end{aligned}$$

(b) Show that the expression of the variance found in (a) in the case of simple linear regression ($p = 1$) can be written as

$$\text{Var} \hat{\epsilon}_{n+1} = \sigma^2 \left(1 + \frac{1}{n} + \frac{(x_{n+1} - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right).$$

(c) For which value of x_{n+1} is the variance of the prediction error minimum? What is the value of the variance in that case?

Our goal is to generalise the result found in (c) in the multiple linear regression setting. We adopt the following notation

$$\mathbf{X} = \begin{pmatrix} 1 & \mathbf{z}_1^t \\ \vdots & \vdots \\ 1 & \mathbf{z}_n^t \end{pmatrix} = (1 \quad \mathbf{Z}_1 \quad \dots \quad \mathbf{Z}_p) = (1 \quad \mathbf{Z}),$$

where the \mathbf{Z}_j are column vectors and \mathbf{Z} is a $n \times p$ matrix. The column means of \mathbf{Z} are put into a vector $\bar{\mathbf{x}}^t = (\bar{x}_1, \dots, \bar{x}_p)$.

(d) Express $\mathbf{X}^t \mathbf{X}$ as a 2×2 block matrix, in terms of \mathbf{Z} , $\bar{\mathbf{x}}$ and n .

We recall the inversion formula for block matrices. Let \mathbf{M} be an invertible matrix, such that

$$\mathbf{M} = \begin{pmatrix} \mathbf{T} & \mathbf{U} \\ \mathbf{V} & \mathbf{W} \end{pmatrix},$$

with \mathbf{T} invertible. Then $\mathbf{Q} = \mathbf{W} - \mathbf{V}\mathbf{T}^{-1}\mathbf{U}$ is invertible and

$$\mathbf{M}^{-1} = \begin{pmatrix} \mathbf{T}^{-1} + \mathbf{T}^{-1}\mathbf{U}\mathbf{Q}^{-1}\mathbf{V}\mathbf{T}^{-1} & -\mathbf{T}^{-1}\mathbf{U}\mathbf{Q}^{-1} \\ -\mathbf{Q}^{-1}\mathbf{V}\mathbf{T}^{-1} & \mathbf{Q}^{-1} \end{pmatrix}.$$

(e) Express $(\mathbf{X}^t \mathbf{X})^{-1}$ as a 2×2 block matrix, in terms of n , $\bar{\mathbf{x}}$ and some matrix Γ^{-1} . Give an expression of Γ in terms of \mathbf{Z} , $\bar{\mathbf{x}}$ and n .

(f) Let $\mathbf{x}_{n+1}^t = (1 \quad \mathbf{z}_{n+1}^t)$ be a new observation. Show that the variance of the prediction error is

$$\text{Var } \hat{\epsilon}_{n+1} = \sigma^2 \left(1 + \frac{1}{n} + \frac{1}{n} (\mathbf{z}_{n+1} - \bar{\mathbf{x}})^t \Gamma^{-1} (\mathbf{z}_{n+1} - \bar{\mathbf{x}}) \right).$$

(g) Assume that Γ is symmetric positive definite. For which value of \mathbf{x}_{n+1} is the variance of the prediction error minimal? What is the value of the variance in that case?

(h) Show that if $\mathbf{X}^t \mathbf{X}$ is invertible, then Γ is indeed symmetric positive definite.

Problem 10.

We consider n observations y_1, \dots, y_n and their associated k -dimensional input vectors $\mathbf{x}_i = (x_{i1}, \dots, x_{ik})^t$, for $i = 1, \dots, n$. We assume that for all i , y_i corresponds to the observed value of a random variable Y_i , and that there exists β such that

$$Y_i \sim \mathcal{N}(\mathbf{x}_i^t \beta, \sigma_i^2), \quad 1 \leq i \leq n,$$

where

- $\beta = (\beta_1, \dots, \beta_k)^t \in \mathbb{R}^k$,
- the Y_i are mutually independent.

Finally, the values σ_i^2 depend on which of the p sub-populations the variables belong to. Put

- $I_1 = \{1, \dots, n_1\}$, the indices of the n_1 elements of the first sub-population,
- $I_2 = \{n_1 + 1, \dots, n_1 + n_2\}$, the indices the n_2 elements of the second sub-population,
- ...
- $I_l = \{n_1 + \dots + n_{l-1} + 1, \dots, n_1 + \dots + n_{l-1} + n_l\}$, the indices of the n_l elements of the l -th sub-population,
- ...
- $I_p = \{n_1 + \dots + n_{p-1} + 1, \dots, n\}$, the indices of the n_p elements of the last sub-population.

We make the following hypothesis: if $i \in I_l$, then $\sigma_i^2 = l\sigma^2$. In other words, the n_1 variables belonging to the first population have variance σ^2 , the n_2 variables belonging to the second population have variance $2\sigma^2$, and so on.

Our goal is to estimate β and σ^2 using maximum likelihood. We denote by $\hat{\beta}$ and $\hat{\sigma}^2$ these estimators.

- (a) Write down the density $f_{Y_i}(y_i)$ of the variable Y_i .
- (b) Show that $\hat{\beta}$ and $\hat{\sigma}^2$ solve the following system of equations

$$\sum_{l=1}^p \frac{1}{l} \sum_{i \in I_l} (y_i - \mathbf{x}_i^t \beta)^2 = n\sigma^2,$$

$$\sum_{l=1}^p \frac{1}{l} \sum_{i \in I_l} (y_i - \mathbf{x}_i^t \beta) x_{ij} = 0, \quad \forall j = 1, \dots, k.$$

- (c) Show that the system derived in (b) can be put into a matrix form,

$$\|A(Y - X\beta)\|^2 = n\sigma^2,$$

$$X^t A^2 (Y - X\beta) = 0,$$

where $\|\cdot\|^2$ denotes the usual Euclidean norm in \mathbb{R}^n , X is the $n \times k$ matrix of observations, Y is the $n \times 1$ vector of output values, and A is an $n \times n$ matrix. Give an expression for A .

- (d) Assuming that $(X^t A^2 X)$ is invertible, give an expression for $\hat{\beta}$ and $\hat{\sigma}^2$.
- (e) Show that $n\hat{\sigma}^2 = \|V\|^2$, where V follows a centered normal distribution.
- (f) Show that $\mathbf{E}\|V\|^2$ is equal to the trace of the covariance matrix of V .
- (g) Show that $n\hat{\sigma}^2/(n - k)$ is an unbiased estimator of σ^2 .
- (h) We denote by X_l the $n_l \times k$ matrix corresponding to the rows of indices I_l of X , assumed full rank, and Y_l the $n_l \times 1$ vector of components of indices I_l of Y . Put $\hat{\beta}_l = (X_l^t X_l)^{-1} X_l^t Y_l$. Show that $\hat{\beta}_l$ is an unbiased estimator of β .

(i) What can we say about the difference of the covariance matrices of $\hat{\beta}_i$ and $\hat{\beta}$?

Hint: Recall that the matrix inequality $B \preceq C$ holds if $C - B$ is positive semidefinite. If B and C are two symmetric positive semidefinite matrices such that $B \preceq C$, then $C^{-1} \preceq B^{-1}$.

Problem 11.

The *UK Building Research Station* collected data on weekly gaz consumption and average external mean temperature in a district in South-East England over a few months. A linear regression explaining gaz consumption as a function of the temperature is carried out in R:

```
Residuals:
      Min       1Q   Median       3Q      Max
-0.97802 -0.11082  0.02672  0.25294  0.63803

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  4.72385     0.12974      ?    < 2e-16 ***
Temp        -0.27793         ?    -11.04 1.05e-11 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3548 on 28 degrees of freedom
Multiple R-Squared:  0.8131,    Adjusted R-squared:  0.8064
F-statistic: 121.8 on 1 and 28 DF,  p-value: 1.046e-11
```

- (i) Write down the model and its assumptions.
- (ii) Fill in the blanks in the R output.
- (iii) Let $Y \sim t_{28}$. What is $\mathbf{P}(|Y| > 11.04)$?
- (iv) Describe the test associated with the row Temp in the R output (the null hypothesis, the alternative, the law under the null, the decision rule).
- (v) Give an interpretation for Multiple R-squared: 0.8131.
- (vi) Give an estimate of the variance of the error term in the simple linear regression model.
- (vii) Explain and interpret the last line
F-statistic: 121.8 on 1 and 28 DF, p-value: 1.046e-11
- (viii) Do you believe that the outside temperature has an effect of gaz consumption? Justify your answer.

Problem 12.

We are interested in the model $Y = X\beta + \epsilon$ under usual conditions. We obtained the following fit, on a learning sample of size $n = 21$:

$$\hat{y} = 6.683_{(2.67)} + 0.44_{(2.32)}x_1 + 0.425_{(2.47)}x_2 + 0.171_{(2.09)}x_3 + 0.009_{(2.24)}x_4,$$

and $R^2 = 0.54$. For each coefficient, the number between brackets represents the absolute value of the test statistic.

- (i) What are the assumptions made on the model?
- (ii) Test for $\beta_1 = 0$ at 5% level.
- (iii) Can you test for $\beta_3 = 1$ against the two-sided alternative $\beta_3 \neq 1$?
- (iv) Test for $\beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$ at 5% level.